



## Calhoun: The NPS Institutional Archive

---

Faculty and Researcher Publications

Faculty and Researcher Publications Collection

---

2015-12-16

# GPU Accelerated Spectral Element Methods: 3D Euler equations

Giraldo, Francis

---

GPU Accelerated Spectral Element Methods: 3D Euler equations American Geophysical Union (AGU) Fall Meeting, 16 December 2015  
<http://hdl.handle.net/10945/48821>



Calhoun is a project of the Dudley Knox Library at NPS, furthering the precepts and goals of open government and government transparency. All information contained herein has been approved for release by the NPS Public Affairs Officer.

**Dudley Knox Library / Naval Postgraduate School**  
**411 Dyer Road / 1 University Circle**  
**Monterey, California USA 93943**

<http://www.nps.edu/library>

# GPU Accelerated Spectral Element Methods: 3D Euler equations

Daniel Abdi, Lucas Wilcox, Timothy Warburton\* & Francis Giraldo  
Naval Postgraduate School, Virginia Tech University\*

**Contact Information:**  
Department of Applied Mathematics  
Naval Postgraduate School  
1 University Circle, Monterey, CA, USA  
Phone: +1 (831) 582 7010  
Email: dsabdi@nps.edu



## Abstract

A GPU accelerated model discontinuous Galerkin method for the solution of the 3D Euler equations in the Non-hydrostatic Unified Model of the Atmosphere (NUMd) is presented. We use algorithms suitable for the single instruction multiple thread architecture of GPUs to accelerate the dynamical core of NUMd by two orders of magnitude relative to one core of a CPU. Tests on one node of the Titan supercomputer yield a speedup of upto 15X on one K20X GPU relative to that on an AMD Opteron CPU with 16 cores. The scalability of the multi-GPU implementation is tested using 16384 GPUs, which resulted in a weak scaling efficiency of about 90%. For portability to heterogeneous computing environment, we used a new programming language OCCA, which can be cross-compiled to either **OpenCL**, **CUDA** or **OpenMP** at runtime. Finally, the accuracy and performance of our GPU implementations are verified using benchmark problems representative of different scales of atmospheric dynamics.

## Introduction

The National Oceanic and Atmospheric Administration (NOAA) has set a goal of 3 km resolution for an operational NWP model within the next decade [4]. The Non-hydrostatic Unified Model of the Atmosphere (NUMd) has recently achieved this goal using 3 million MPI threads of the *multi-core* MIRA supercomputer [3]. The current study aims to port NUMd to *many-core* architecture (CPUs, Xeophi etc.) using a new programming language called OCCA [12] that provides backends to different kinds of hardware. This will enable NUMd to use future heterogeneous hardware using the same code that will be compiled and optimized for each computing unit in the heterogeneous cluster.

NUMd uses element based Galerkin (EBG) methods, continuous Galerkin (CG) and discontinuous Galerkin (DG), for spatial discretization. These methods are well suited for GPU computing for two reasons a) Localized memory accesses result in low communication overhead. In contrast, global methods require an *all-to-all* communication that severely degrades scalability. b) High order polynomial expansion of the solution results in large arithmetic intensity per degree of freedom. Parallel implementation of DG is often easier and more efficient than that of CG because of a smaller communication stencil; only neighbors sharing a face need to communicate in DG as opposed to edge and corner neighbor communication required by CG. Moreover, DG allows for a simple overlap of computation of volume integrals and intra-processor flux with communication of boundary data, which can be exploited to improve the efficiency of the parallel implementation [1].

## Governing equations

The Euler equations are

$$\begin{aligned} \frac{\partial \rho}{\partial t} + \nabla \cdot (\rho \mathbf{u}) &= 0 \\ \frac{\partial \rho \mathbf{u}}{\partial t} + \nabla \cdot (\rho \mathbf{u} \otimes \mathbf{u} + P \mathbf{I}) &= -\rho \mathbf{g} \\ \frac{\partial \rho h}{\partial t} + \nabla \cdot (\rho h \mathbf{u}) &= 0 \end{aligned}$$

or in compact vector notation form

$$\frac{\partial \mathbf{q}}{\partial t} + \nabla \cdot \mathbf{F} = \mathbf{S}(\mathbf{q})$$

with the equation of state given by

$$P = P_0 \left( \frac{\rho h}{P_0} \right)^\gamma$$

## Discretization

Continuous Galerkin:

$$\int_{\Omega_e} \psi_i \frac{\partial \mathbf{q}_N}{\partial t} d\Omega_e + \int_{\Omega_e} \psi_i \nabla \cdot \mathbf{F} d\Omega_e = \int_{\Omega_e} \psi_i \mathbf{S}(\mathbf{q}_N) d\Omega_e$$

Discontinuous Galerkin:

$$\int_{\Omega_e} \psi_i \frac{\partial \mathbf{q}_N}{\partial t} d\Omega_e + \int_{\Gamma_e} \psi_i \hat{\mathbf{n}} \cdot (\mathbf{F}^* - \mathbf{F}) d\Gamma_e + \int_{\Omega_e} \psi_i \nabla \cdot \mathbf{F} d\Omega_e = \int_{\Omega_e} \psi_i \mathbf{S}(\mathbf{q}_N) d\Omega_e$$

The major compute kernels are

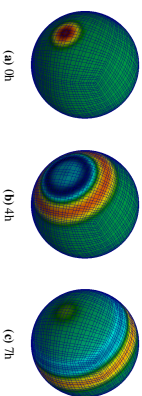
1. Volume integration kernel
2. Surface integration kernel
3. Time update kernel
4. Direct Stiffness Summation (DSS) required by CG

## Speedup

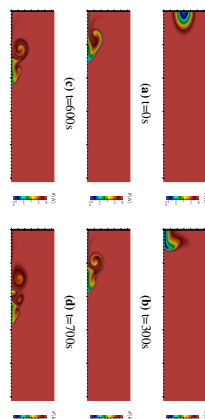
**Table 1:** Speedup comparison between CPU and GPU for single and double precision calculations at different number of elements and polynomial orders using CUDA translation of OCCA kernel code. A maximum speedup of about 15X is observed using the K20X GPU card on the Titan supercomputer relative to a 16-core 2.2GHz AMD Opteron 6274 CPU. The CPU/GPU times (seconds) and speedups are given first for the double precision result and then for single precision.

N	30x30x900 elements			40x40x1600 elements		
	CPU	Speedup	GPU	CPU	Speedup	GPU
2	10.62/0.98	2.17/1.57	4.90/6.36	18.8/17.41	3.70/2.53	5.09/6.88
3	22.01/1.966	3.06/1.87	7.19/10.51	41.5/34.72	5.04/5.06	8.23/11.43
4	46.4/3.819	5.12/3.03	9.07/11.61	81.9/103.55	8.69/5.07	10.47/12.85
5	77.03/6.135	9.88/4.86	7.80/12.62	137.49/107.30	17.1/18.33	14.83/17.25
6	122.27/9.567	16.11/7.71	7.59/12.41	210.35/166.40	28.15/13.49	17.47/12.33
7	195.61/13.521	18.87/9.65	10.37/14.01	343.74/236.09	33.05/16.86	18.00/14.06
N	80x80=6400 elements			120x120=14400 elements		
	CPU	Speedup	GPU	CPU	Speedup	GPU
2	80.72/7.041	13.3/8.94	6.05/7.38	17.82/19.78	6.02/6.78	9.80/10.91
3	179.07/14.219	18.46/11.18	9.70/12.72	41.08/24.87	9.86/10.91	13.97/14.54
4	330.54/26.68	32.71/19.02	10.71/14.13	72.77/42.34	10.89/11.93	14.83/15.75
5	587.17/42.66	67.03/23.98	8.76/13.27	132.97/100.73	15.05/72.08	18.36/19.49
6	925.25/66.25	110.92/52.91	8.34/13.16	208.68/158.65	24.95/118.39	26.74/17.49
7	1406.61/96.81	130.16/66.41	10.81/14.58	315.84/222.72	29.30/148.76	30.76/19.49

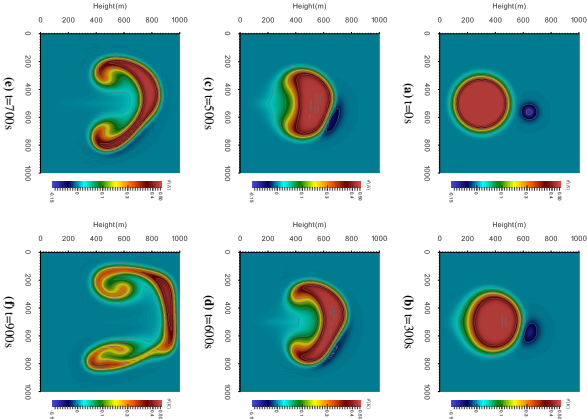
## Verification



**Figure 1:** Propagation of an acoustic wave. The density perturbation after 0 hour, 4 hours and 7 hours. A cubed sphere grid with 10x10 elements in the horizontal and 3 elements in the vertical with 3rd degree polynomial is used.

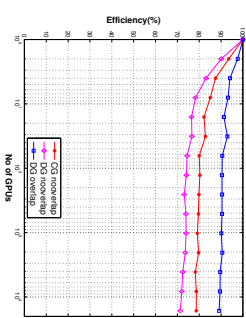


**Figure 2:** Density current. Evolution of potential temperature perturbation  $\theta'(K)$  from with CG and an artificial viscosity of  $\nu = 75 m^2/s$  for stabilization. Results are shown at t=300, 600, 700, 800 and 900 seconds. A grid of 128x132 elements with 7 degree polynomials is used for an effective resolution of 50m in x and z directions.



**Figure 3:** Colliding thermal bubbles. Evolution of potential temperature perturbation  $\theta'(K)$  from with CG and an artificial viscosity of  $\nu = 1.5 m^2/s$  for stabilization. Results are shown at t=300, 500, 600, 700 and 900 seconds. A grid of 10x10 elements with 6th degree polynomials is used.

## Scalability



**Figure 4:** Scalability test of Multi-GPU implementation on Titan. The scalability of NUMd for up to 16384 GPUs on the Titan supercomputer is shown. Efficiency of Titan supercomputer is shown. Efficiency of about 90% is observed for CG and DG. The efficiency of DG is significantly improved by about 20% when overlapping communication with computation, which helps to hide both latency in the data copy between the CPU and GPU and the CPU-GPU communication.

## Conclusions

- NUMd is ported to many-core architecture using a new programming language called OCCA that can be cross compiled to CUDA, OpenCL, OpenMP
- A speed up of upto 240X relative to a single core CPU (or 15X relative to a 16-core CPU) is obtained using a single NVIDIA K20X GPU
- A weak scaling efficiency of about 90% is observed using 16384 GPUs of the Titan supercomputer

## References

- [1] J. F. Kelly and F. X. Giraldo. Continuous and discontinuous galerkin methods for a scalable three-dimensional nonhydrostatic atmospheric model: limited area mode. *J. Comput. Phys.*, 231:7988–8008, 2012.
- [2] D. Medina, St-Cyr Anilk, and T. Warburton. OCCA: A unified approach to multi-threading languages. *arXiv*, 2014.
- [3] A. Müller, M. Köpcke, S. Marras, L.C. Wilcox, T. Isaac, and F.X. Giraldo. Strong scaling for numerical weather prediction at petascale with the atmospheric model numd. *Submitted to : 30th IEEE International Parallel and Distributed Processing Symposium*, 2016.
- [4] NOAA. High performance computing plan 2015-2020. *National Oceanic and Atmospheric Administration*, pages 1–11, 2015.

## Acknowledgements

The authors gratefully acknowledge the support of the Office of Naval Research through program element PE-0602435N. This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725.